

# Mini-batch gradient descent

→ If mini-batch size =  $m \Rightarrow$  Batch Gradient Descent  $(X^{(1)}, y^{(1)}) = (X, y)$

" " = 1  $\Rightarrow$  Stochastic Gradient Descent.

Every example is a mini-batch.  $(X^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)})$

In practice,  $m$  in mini-batch gradient descent is somewhat in between 1 and  $m$

Stochastic  
gradient  
descent



lose speedup  
from Vectorization.

In-between

(mini-batch size not  
too big or small)



Fastest Learning

- Vectorization ( $\sim 1000$ )
- Make progress without  
processing entire training set.

Batch

gradient descent  
(mini-batch size =  $m$ )



Too long

per iteration

Mini-Batch gradient descent:

Repeat  $n$  epochs

for  $t = 1, \dots, 5000$ : mini-batch iterations.

forward prop on  $X^{(t)}$

Compute Cost  $J^{(t)} = \frac{1}{1000} \sum_{i=1}^N L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2000} \sum \|W^{(2)}\|_F^2$

backprop to compute gradients

$W^{[2]} := W^{[2]} - \alpha dw$ ,  $b^{[2]} := b^{[2]} - \alpha db$

Choosing mini-batch size:

if small training set  $\rightarrow$  use batch gradient descent  
( $m < 2000$ ).

Typical mini-batch Size

$\rightarrow 64, 128, 256, 512$

$2^6 \quad 2^7 \quad 2^8 \quad 2^9$

Theoretically  $2^n$  makes gradient descent run faster.